

# PATENT ABSTRACTS OF JAPAN

(11)Publication number : 11-119791

(43)Date of publication of application : 30.04.1999

(51)Int.Cl.

G10L 3/00  
G10L 3/00  
G10L 3/00  
// A63F 9/22  
G06F 17/28

(21)Application number : 09-286372

(71)Applicant : HITACHI LTD  
HITACHI ULSI SYSTEMS CO LTD

(22)Date of filing : 20.10.1997

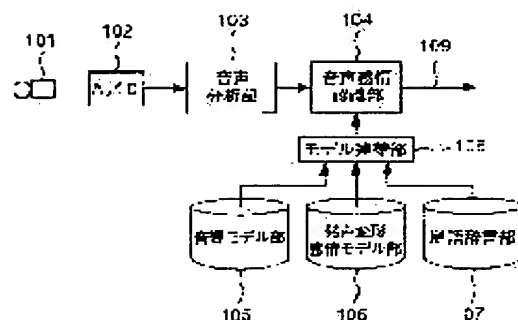
(72)Inventor : WAKIZAKA SHINJI  
KONDO KAZUO  
OBUCHI YASUNARI  
TOUSHITA TETSUJI  
ISHIKAWA YASUYO

(54) SYSTEM AND METHOD FOR VOICE FEELING RECOGNITION

(57)Abstract:

PROBLEM TO BE SOLVED: To recognize the level of a speaker's feeling by a voice recognition system.

SOLUTION: The system and method are equipped with a dictionary part 107 where object words of voice recognition are gathered, a voice analysis part 103 which performs a voice analyzing process, a sound model part 105 which has patterns of voice in phoneme units, a voicing deformation feeling model part 106 which represents the deformation of a vocal sound spectrum by a feeling, and a voice recognition part 104 which performs a voice recognizing process by coupling the sound model part 105, voicing deformation feeling model part 106, and dictionary part 107, and outputs an object word of voice recognition as a voice recognition result and also outputs a feeling level representing the degree of the speaker's feeling that the voice has. Other voice analysis parts output feeling levels from the features of the power of the voice.



## LEGAL STATUS

[Date of request for examination]

[Date of sending the examiner's decision of rejection]

[Kind of final disposal of application other than the examiner's decision of rejection or application converted registration]

[Date of final disposal for application]

[Patent number]

[Date of registration]

[Number of appeal against examiner's decision of rejection]

[Date of requesting appeal against examiner's decision  
of rejection]

[Date of extinction of right]

Copyright (C); 1998,2000 Japan Patent Office

(19) 日本国特許庁 (J P)

(12) 公開特許公報 (A)

(11) 特許出願公開番号

特開平11-119791

(43) 公開日 平成11年(1999) 4月30日

(51) Int.Cl.<sup>6</sup>  
G 1 0 L 3/00

識別記号  
5 3 1

F I  
G 1 0 L 3/00

5 3 1 N

R

5 3 5

5 3 5

// A 6 3 F 9/22  
G 0 6 F 17/28

A 6 3 F 9/22  
G 0 6 F 15/38

F

V

審査請求 未請求 請求項の数 8 O L (全 13 頁)

(21) 出願番号 特願平9-286372

(22) 出願日 平成 9 年(1997) 10月20日

(71) 出願人 000005108

株式会社日立製作所

東京都千代田区神田駿河台四丁目 6 番地

(71) 出願人 000233169

株式会社日立超エル・エス・アイ・システムズ

東京都小平市上水本町 5 丁目22番 1 号

(72) 発明者 脇坂 新路

東京都小平市上水本町五丁目20番 1 号株式会社日立製作所半導体事業部内

(74) 代理人 弁理士 高橋 明夫 (外 1 名)

最終頁に続く

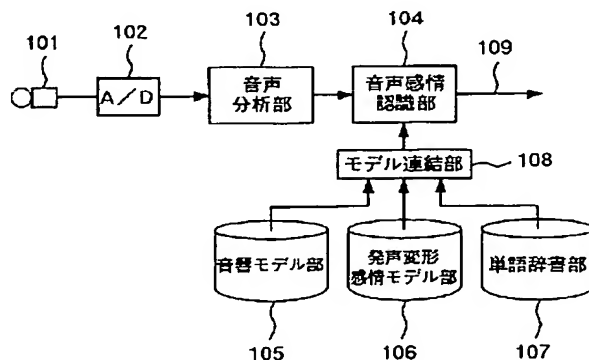
(54) 【発明の名称】 音声感情認識システムおよび方法

(57) 【要約】

【課題】 音声認識システムにおいて話者の感情のレベルを認識する。

【解決手段】 音声認識の対象となる単語を集めた辞書部と、音声分析処理を行う音声分析部と、音声のパターンを音素単位でもつ音響モデル部と、感情による音韻スペクトルの変形を表す発声変形感情モデル部と、音声分析結果に対して、音響モデル部と発声変形感情モデル部と辞書部とを連結して音声認識処理を行う音声認識部とを備え、音声の特徴から、音声認識の対象となる単語を音声認識結果として出力すると共に、音声もっている話者の感情の度合を示す感情レベルを出力する。他の音声分析部は音声のパワーの特徴から感情レベルを出力する。

図 1



## 【特許請求の範囲】

【請求項1】音声認識の対象となる単語や文章を集めて辞書として定義し、音声認識結果として、それらの単語や文章を辞書部からピックアップして、文字列表示や音声合成を用いて出力する音声認識システムにおいて、取り込んだ音声に対して音声分析処理を行う音声分析部と、音声のパターンを音素単位でもつ音響モデル部と、感情による音韻スペクトルの変形を表す発声変形感情モデル部と、音声分析結果に対して、音響モデル部と発声変形感情モデル部と辞書部とを連結して、音声認識処理を行う音声認識部とを備え、音声の特徴から、音声認識の対象となる単語や文章を音声認識結果として出力するとともに、音声もっている話者の感情の度合を示すレベルを出力することを特徴とする音声感情認識システム。

【請求項2】請求項1記載の音声感情認識システムにおいて、音声もっている話者の感情の度合を示すレベルは、数字0～N（Nは整数）であることを特徴とする音声感情認識システム。

【請求項3】請求項1記載の音声感情認識システムにおいて、前記辞書部は、音声認識の対象となる単語や文章を集めた辞書と、それらの単語や文章に対して、音声もっている感情のレベルを表現する修飾語を集めた辞書とを備え、音声認識結果として、それらの単語や文章をピックアップするとともに、感情のレベルを表現する修飾語をピックアップして、単語や文章に修飾語を付加して、文字や音声合成を用いて出力することを特徴とする音声感情認識システム。

【請求項4】音声認識の対象となる単語や文章を集めて辞書として定義し、音声認識結果として、それらの単語や文章をピックアップして、文字列表示や音声合成を用いて出力する音声認識システムにおいて、取り込んだ音声に対して音声分析処理を行う音声分析部と、音声のパターンを音素単位でもつ音響モデル部と、音声分析結果に対して音響モデル部と辞書部とを連結して音声認識処理を行う音声認識部とを備え、取り込んだ音声に対して、音声分析処理を行う音声分析部は、感情の度合が現われる音の強弱を示すパワーの特徴から、感情の度合を示すレベルを出力することを特徴とする音声感情認識システム。

【請求項5】音声認識の対象となる単語や文章を集めた辞書と、取り込んだ音声に対して音声分析処理を行う音声分析部と、音声のパターンを音素単位でもつ音響モデル部と、感情による音韻スペクトルの変形を表す発声変形感情モデル部と、音声分析結果に対して、音響モデル部と発声変形感情モデル部と辞書部とを連結して音声認識処理を行う音声認識部とを備え、音声の特徴から音声認識の対象となる単語や文章を音声認識結果として出力

するとともに、発声変形感情モデル部からのデータを用いて音声もっている話者の感情の度合を出力することを特徴とする音声感情認識方法。

【請求項6】音声認識の対象となる単語や文章を集めて辞書と、取り込んだ音声に対して音声分析処理を行うと共に音のパワーを分析して感情のレベルを出力することが出来る音声分析部と、音声のパターンを音素単位でもつ音響モデル部と、音声分析結果に対して音響モデル部と辞書を連結して音声認識処理を行う音声認識部とを備え、音声分析部は取り込んだ音声に対して感情の度合が現われる音の強弱を示すパワーの特徴から、感情の度合を認識し、この出力感情のレベルを出力することを特徴とする音声感情認識方法。

【請求項7】請求項5又は6記載の音声感情認識方法において、音声もっている話者の感情の度合を示すレベルは、数字0～N（Nは整数）であることを特徴とする音声感情認識方法。

【請求項8】請求項5又は6記載の音声感情認識方法に於いて、前記辞書部は、音声認識の対象となる単語や文章を集めた辞書と、それらの単語や文章に対して、音声もっている感情のレベルを表現する修飾語を集めた辞書とを備え、音声認識結果として、それらの単語や文章をピックアップするとともに、感情のレベルを表現する修飾語をピックアップして、単語や文章に修飾語を付加して、文字や音声合成を用いて出力することを特徴とする音声感情認識システム。

## 【発明の詳細な説明】

## 【0001】

【発明の属する技術分野】本発明は、音声認識システムおよび方法に係わり、カーナビゲーションシステム、車載用PC、PDA（パーソナル・デジタル・アシスタント）、ハンドヘルドPCに代表される小型情報機器、携帯型音声翻訳機、ゲーム、家電機器に用いる音声認識システムであって、特に、音声認識の対象となる単語や文章の認識とともに、感情を表わす単語や文章においては、感情の度合まで認識する音声感情認識システムおよび方法に関する。

## 【0002】

【従来の技術】近年、音声認識技術を用いた小型情報システムが普及しつつある。カーナビゲーションシステムをはじめとして、PDAに代表される小型情報機器、携帯型翻訳機等である。

【0003】このような音声認識システムの例として、特願平5-35776号公報の「言語自動選択機能付翻訳装置」には、マイクから入力した操作者の音声を認識して、翻訳し、翻訳した言語の音声を出力するようにした携帯用の翻訳装置に関する技術が開示されている。

【0004】以下、図7を用いてこのような従来技術に

係わる音声翻訳装置の概要について説明する。

【0005】図7は従来技術に係わる音声翻訳装置の構成を示すブロック図である。制御部701は、マイクロプロセッサ等からなり、装置の各部を制御する。音声区間切出し部702は、マイク709から入力された音声デジタル信号に変換して切り出し、音声認識部703に送る。音声認識部703は、キーボード又はスイッチ等による操作信号711を受けた制御部701の指示により、マイク709、音声区間切出し部702を経て、切り出された音声进行分析する。そしてその結果を、音声認識辞書部707に格納された標準音声パターンと比較することにより、音声認識をおこなう。音声合成部705は、音声認識部703により認識された音声に対応した翻訳語を、翻訳語データ用メモ리카ード706から読み込み、音声信号に変換してスピーカアンプ710、スピーカ708を経て出力する。

【0006】表示部704は、翻訳装置の使用者への指示や翻訳語の文字による表示等をおこなう。翻訳語データ用メモ리카ード706は、ROMカード等からなり、翻訳語を音声合成して出力する場合には、音声データを格納している。また、この翻訳語データ用メモ리카ード706から、翻訳語に対応したキャラクターコードを読み込み、表示部704に表示する。そして、この翻訳語データ用メモ리카ード706を他の言語のものと交換することにより、複数の言語に翻訳することが可能となる。音声認識辞書部707は、RAM等からなり、操作者の発声に応じた標準音声パターンを格納している。この標準音声パターンは、操作者があらかじめ格納しておく。

【0007】

【発明が解決しようとする課題】このような音声認識、音声合成技術の分野は、半導体技術の向上を背景として、システムがより人間的なユーザインタフェースを提供すべきであるという要望から、その発展が期待されている。上記従来の音声認識技術を用いた小型情報システムにおいても、カーナビゲーションシステムをはじめとして、PDAに代表される携帯型情報機器、携帯型翻訳機、さらに、音声インタフェースを持った情報家電として、今後ますます普及してくることが予想される。

【0008】しかしながら、音声認識は、処理すべき情報量が膨大なものになるため、従来の技術では、認識率や認識応答時間の性能を低下させないためには、認識する語数に制約を設ける必要がある。そのためには、あらかじめ登録しておいた単語、文に対して、その文字列が持つ統計的な話者の音声の特徴と、実際に話者が発声した音声の特徴とを比較し、確率的に一番近い値を認識結果としている。

【0009】今後、音声認識における技術革新や、それを実現するソフトウェア、ハードウェアの性能向上により、認識率や認識応答時間の性能は向上することが考え

られる。そこで、さらに、人間的なユーザインタフェースを提供するためには、単に、従来の音声認識技術において、あらかじめ登録した単語、文の文字列を認識するだけでなく、話者の感情や意図を認識できれば、たとえ制限された認識語数においても、使い勝手の向上が期待できる。しかしながら、従来の音声認識システムでは、あらかじめ登録した単語や文の文字列のみを音声で照合して、入力した音声に最も近い文字列を音声認識結果として出力する音声認識システムであり、音声を発声した話者の感情や意図までは認識できない。

【0010】本発明は、システムが少しでも人間的なユーザインタフェースを持てるように、上記問題点を解決するためになされたものである。

【0011】本発明の目的は、小型情報システムに用いられる音声認識システムにおいて、入力された音声に対して、辞書に登録された単語や文の文字列を認識するとともに、入力された音声を持っている話者の感情や意図を認識することができる音声感情認識システム及び方法を提供することにある。

【0012】また、本発明の他の目的は、小型情報システムに用いられる音声認識システムにおいて、入力された音声を持っている話者の感情や意図を感情の度合いを表現する数字や修飾語に変換して、人間とシステムにおける良好な音声インタフェースを実現することである。

【0013】

【課題を解決するための手段】上記目的を達成するために、本発明の音声感情認識システム及び方法に於いては音声認識の対象となる単語や文章を集めて辞書として定義し、音声認識結果として、それらの単語や文章を辞書部からピックアップして、文字列表示や音声合成を用いて出力する音声認識システムにおいて、取り込んだ音声に対して音声分析処理を行う音声分析部と、音声のパターンを音素単位でもつ音響モデル部と、感情による音韻スペクトルの変形を表す発声変形感情モデル部と、音声分析結果に対して音響モデル部と発声変形感情モデル部と辞書部とを連結して音声認識処理を行う音声認識部とを備え、音声の特徴から音声認識の対象となる単語や文章を音声認識結果として出力するとともに、音声を持っている話者の感情の度合を出力するようにしたものである。

【0014】より詳しい1実施例に於いては、音声を持っている話者の感情の度合を示すレベルは、数字0～N（Nは整数）であるようにしたものである。

【0015】また、本発明の音声感情認識システム及び方法に於いては、音声を持っている話者の感情の度合は、音声認識の対象となる単語や文章を集めた辞書とそれらの単語や文章に対して感情のレベルを表現する修飾語を集めた辞書と有する辞書部を備え、音声認識結果として、それらの単語や文章をピックアップするとともに、感情のレベルを表現する修飾語をピックアップし

て、単語や文章に修飾語を付加して、文字や音声合成を用いて出力するようにしたものである。

【0016】さらに詳しい1実施例に於いては、音声認識の対象となる単語や文章を集めて辞書として定義し、音声認識結果として、それらの単語や文章をピックアップして、文字列表示や音声合成を用いて出力する音声認識システムにおいて、取り込んだ音声に対して音声分析処理を行う音声分析部と、音声のパターンを音素単位でもつ音響モデル部と、音声分析結果に対して音響モデル部と辞書部とを連結して音声認識処理を行う音声認識部とを備え、取り込んだ音声に対して音声分析処理を行う音声分析部は、感情の度合いが現われる音の強弱を示すパワーの特徴から感情の度合いを出力することが出来る。

【0017】

【発明の実施の形態】以下、本発明に係る各実施形態を図1から図6を用いて説明する。

【0018】図1は本発明に係る音声および感情認識システムの各機能とその処理の流れを示すブロック図である。

【0019】音声および感情認識をおこなうために、図1に示されるマイク101から音声を取り込まれる。取り込まれた音声であるアナログ信号は、アナログ信号をデジタル信号に変換するA/D変換器102によって、任意に決められたサンプリング周期により、アナログデータからデジタルデータに変換される。変換された音声のデジタルデータは、音声分析部103によって、雑音処理や音声分析や話者適応などの前処理がなされ、音声感情認識部104により音声および感情認識がなされる。ここで、音声および感情認識とは、2つの処理を実行する。

【0020】第1の処理は、音声信号を解析して、それを短い時間ごとの音素として分析して、そのパターンを解析し、該当する単語や文章を辞書から選択することである。

【0021】第2の処理は、音声信号を解析して、それを短い時間(5~20ms)ごとの音素として分析して、そのパターンを解析し、話者が発声した音声の感情の度合いを示すレベルを単語や文章ごとに選択することである。

【0022】以上の2つの処理から、音声感情認識システムの出力として、音声認識結果および音声の感情レベル109を生成する。

【0023】音声感情認識部104は、音声分析部103で分析された入力音声の音声分析結果に対して、音響モデル部105、発声変形感情モデル部106、単語辞書部107をモデル連結部108によって連結された音素単位で照合して、単語辞書部107に登録した単語の中で、一番近い単語をピックアップする。さらに、ピックアップされた単語の入力音声を持っている感情の度合いを示すレベルを選択する。なを、図1に示す実施例に

於いては、電源を投入すると、モデル連結部108で連結された音素単位の単語及び感情の度合いを示すレベルは音声感情認識部104に記憶され、音声分析部103からの音声分析結果と直ちに照合出来るようになっていく。

【0024】音響モデル部105は、音声認識に用いられるモデルであり、具体的には、単語辞書部107に用いられている文字と音素との対応であり、音素の特徴が出現する確率の分布、出現した音素の特徴が次のどの特徴が現れる状態に遷移するかの確率の分布を記憶したものである。音素の特徴が出現する確率の分布について説明する。例えば、「あつい」という音声の「あ」に対して音声スペクトラムは人によって変わるため、「あ」という音素に対して、横軸に音声スペクトラムをとり、縦軸に音素が出現する確率を取ると、音声スペクトラムに対して「あ」と認識される確率が変わることを言う。次に、出現した音素の特徴が次のどの特徴が現れる状態に遷移するかの確率の分布について説明する。例えば、「あ」という音素は「あつい」のように次に「つ」がくる場合もあるし、「あさい」のように次に「さ」がくることもあるし、「あまい」のように次に「ま」に遷移することもある。「あ」が次にどの音素に遷移するかの確率は各音素によって変わる。つまり、ある音素の特徴が次にどの音素の特徴に変化するかの確率は変わるので、この確率の分布を言う。

【0025】音響モデル部105は、あらかじめ声を登録しなくても、誰が話し手でもその声を認識できるいわゆる「不特定話者対応」が、一般的になってきている。このような音響モデルとしては、例えば、隠れマルコフモデル(HMM: Hidden Markov Model)を用いることができる。

【0026】発声変形感情モデル部106は、感情の変化による音韻スペクトルの変形要素に着目して、感情が変化したときの単語辞書部107に用いられている文字と音素との対応である。即ち、感情を込めたときに音素の確率の分布が変わるが、その時の音素の特徴が出現する確率の分布、出現した音素の特徴が次のどの特徴が現れる状態に遷移するかの確率の分布を記憶したものである。この、出現した音素の特徴が次のどの特徴が現れる状態に遷移するかの確率の分布は、例えば、「あつい」という言葉を感情を込めて「あつい」と言っても変化しないが、「あちー」とか「あちい」に変化した場合が変わる。このような発声変形感情モデル部106としては、例えば、隠れマルコフモデル(HMM: Hidden Markov Model)を用いることができる。

【0027】単語辞書部107は、言葉、単語(名詞、動詞等)、文章を集めたものである。例えば、カーナビゲーションシステムにおいては、通り名、地名、建造物名、町名、番地、交差点名、個人住宅(個人名)、電話番号等や、必要最小限の会話に必要な言葉の集合体であ

る。ただし、音声認識感情システムでは、特に、単語の中でも、感情を表現する単語、あるいは、感情が現われる単語で構成された単語の集合体である。より具体的には、話者が発声する「暑い」「寒い」「熱い」「冷たい」「はやく」「おそく」「大きい」「小さい」「赤い」「白い」「高く」「低く」「走れ」「進め」「戻れ」「回れ」「飛べ」等の言葉である。また、名詞等の感情表現でない単語も含まれる。この単語辞書部107は、システムの能力に応じて一つの辞書あたり、例えば10～5000語の単語で構成する。

【0028】以上から、音声感情認識とは、音声を解析して、それを短い時間ごとの音素として分析して、そのパターンを解析し、該当する単語や文章を辞書から選択するとともに、話者が発声した音声の感情の度合いを示すレベルを単語や文章ごとに選択することである。

【0029】なお、図1に示す各処理ブロックは、複数のLSIやメモリで構成されたシステムであっても、半導体素子上に構成された一つないし複数のシステムオンチップであってもよい。また、各処理は、専用LSIや専用ICで処理するハードウェアであっても、DSPやRISCマイコン等のソフトウェアで実現したミドルウェアであってもよい。

【0030】図2は、隠れマルコフモデル(HMM:Hidden Markov Model)による日本語音素のモデル化の例である。

【0031】201、202、203は音素分布の状態を表わしている。話者が発声した音声は、「あつい」であり、発音記号の1例で表わす「atsui」である。説明を簡単にするために、図2(a)において、「a」が201の状態に対応し、「tsu」が202の状態に対応し、「i」が203の状態に対応している。実際の音声認識では、状態をさらに細分化して表わしている。音声は、非定常信号であり、あるときは「a」のスペクトル、あるときは「tsu」のスペクトルという具合に、スペクトルの性質が時々刻々と変化することによって言語情報を伝える。この非定常な信号は、性質の異なる定常信号の音素片の連続とみることができる。この性質の異なる定常信号の音素片の一つ一つが、201～203に示したHMM状態遷移ネットワークの状態に対応している。この状態、すなわち、非定常信号源からの出力として音声のスペクトルが観測される。観測値は、短時間フレーム毎の音声信号のLPC分析結果であっても、ベクトル量子化された符号であってもよい。よって、HMMとは、状態201と状態202の間の状態遷移確率207と、状態201から出力される音声のスペクトルが出力される確率201である。確率201とは、「a」の音素分布の内どの確率値が出力されるかと言う事を示す。即ち、201～203は、音素の分布であり、各状態から出力される音声のスペクトルが出力される確率を示したものであ。204～209は、各状態

が次にどの状態に遷移するかの確率を示したものである。この内、204～206は、例えば、ある音を長く発音したとすると、これを音声分析した場合、ある時間間隔の中では、また同じ音に戻ることを示している。

【0032】つぎに、音声認識に用いるHMMの一例を説明する。

【0033】図2(a)に示す曲線211は、状態201から出力される音声のスペクトルが出力される確率を連続分布で表現したものである。ここで、音声のスペクトルは、音声の特徴パラメータを $i$ 次元としたときの $n$ 番目の特徴パラメータとする。つまり、音声の特徴を表わす表現方法としては何種類もあるが、仮にこの表現方法が $i$ 個あったとすると、その $n$ 番目の表現方法の特徴パラメータを意味する。横軸は、状態201から出力される音声のスペクトルであり、縦軸は、その確率値である。この分布は、平均 $\mu_a$ 、分散 $\sigma_a$ をもつ連続分布である。同様に、図2(b)に示す曲線212は状態202から出力される音声のスペクトルが出力される確率を連続分布で表現したものである。横軸は、状態202から出力される音声のスペクトルであり、縦軸は、その確率値である。この分布は、平均 $\mu_{tsu}$ 、分散 $\sigma_{tsu}$ をもつ連続分布である。

【0034】図2(c)に示す曲線213は、状態203から出力される音声のスペクトルが出力される確率を連続分布で表現したものである。横軸は、状態203から出力される音声のスペクトルであり、縦軸は、その確率値である。この分布は、平均 $\mu_i$ 、分散 $\sigma_i$ をもつ連続分布である。

【0035】ここで、認識対象単語として登録された「atsui」の単語辞書部107に話者が「あつい」と音声を入力する。「あ」の音声に対して、音声分析が行われ、音声の特徴が出力される。例えば、音声の特徴パラメータを $i$ 次元としたときの $n$ 番目の特徴パラメータを使用するものとする、「あ」の特徴 $f_{n1}$ が出力される。このとき、単語辞書「a」において、特徴 $f_{n1}$ の出現する確率が連続分布曲線211から計算され、確率値 $p_{n1}$ が出力される。同様にして、「つ」「い」の音声に対して、音声分析が行われ、音声の特徴が出力される。それぞれ、単語辞書「tsu」において、特徴 $f_{n2}$ の出現する確率が連続分布曲線212から計算され、確率値 $p_{n2}$ が出力される。また、単語辞書「i」において、特徴 $f_{n3}$ の出現する確率が連続分布曲線213から計算され、確率値 $p_{n3}$ が出力される。さらに、音素分布状態201から音素分布状態202間の状態遷移確率においても同様の処理がこなわれ、状態遷移先を状態207に決定している。最終的に、登録された単語辞書「atsui」に対して、音声入力された「あつい」の出現する確率値は $P_{atsui} = p_{n1} + p_{n2} + p_{n3}$ となる。この一連の処理を、登録された単語辞書全てにおいて計算し、確

率値の一番高かったものが、認識結果となる。以上が音声認識の一連の処理である。

【0036】さらに、図2(e)から図2(g)を用いて、音声感情認識における発声変形感情モデルを用いたHMMの一例を説明する。

【0037】発声変形感情モデル部106は、感情の変化による音韻スペクトルの変形要素に着目して、感情が変化したときの単語辞書部107に格納されている単語の文字と音素との対応であり、音素の特徴が出現する確率の分布、出現した音素の特徴が次のどの特徴が現れる状態に遷移するかの確率の分布を記憶したものである。

【0038】曲線211は前に説明したように、状態201から出力される音声のスペクトルが出力される確率を連続分布で表現したものである。ここで、音声のスペクトルは、音声の特徴パラメータを $i$ 次元としたときの $n$ 番目の特徴パラメータとする。横軸は、状態201から出力される音声のスペクトルであり、縦軸は、その確率値である。この分布は、平均 $\mu_a$ 、分散 $\sigma_a$ をもつ連続分布である。このとき、音声の特徴パラメータを $i$ 次元としたときの $n$ 番目の特徴パラメータにおいて、感情の変化による音韻スペクトルの変形が顕著に現れたとする。そこで、話者が通常の感情で発声したときの音声スペクトルの連続分布曲線を211とし、話者が感情をこめて発声したとき、すなわち、感情の変化により変形した時の音声スペクトルの連続分布曲線を214とする。よって、従来の音声認識に用いられてきた音響モデルのHMMに加えて、感情の変化により音韻スペクトルの変形が現われる特徴パラメータだけで構成した確率分布を音声感情認識モデルのHMMとして用意する。曲線214は、音声感情認識モデルにおいて、状態201から出力される感情の変化による音声のスペクトルが出力される確率を連続分布で表現したものである。横軸は、状態201から出力される音声のスペクトルであり、縦軸は、その確率値である。この分布は、平均 $\mu_{a_s}$ 、分散 $\sigma_{a_s}$ をもつ連続分布である。ここで、認識対象単語として登録された「atsui」の単語辞書において、実際に、話者が「あつ」と感情をこめて音声を入力する。「あ」の音声に対して、音声分析が行われ、音声の特徴が出力される。例えば、音声の特徴パラメータを $i$ 次元としたときに、 $n$ 番目の音声の特徴パラメータを採用したとすると、「あ」の特徴 $f_{n1_e}$ が出力される。このとき、単語辞書「a」において、特徴 $f_{n1_e}$ の出現する確率が連続分布曲線214から計算され、確率値 $p_{n1_e}$ が出力される。ここで、この連続分布曲線214に関して、確率値 $p_{n1_e}$ は、話者が通常の発声をしたときの特徴 $f_{n1}$ の確率値 $p_{n1}$ より高い値をとる。

【0039】また、曲線222は、状態202から出力される音声のスペクトルが出力される確率を連続分布で表現したものである。横軸は、状態202から出力され

る音声のスペクトルであり、縦軸は、その確率値である。この分布は、平均 $\mu_{tsu_s}$ 、分散 $\sigma_{tsu_s}$ をもつ連続分布である。曲線223は、状態203から出力される音声のスペクトルが出力される確率を連続分布で表現したものである。横軸は、状態203から出力される音声のスペクトルであり、縦軸は、その確率値である。この分布は、平均 $\mu_{i_s}$ 、分散 $\sigma_{i_s}$ をもつ連続分布である。

【0040】曲線214の場合と同様に、「つ」「い」の音声に対して、音声分析が行われ、音声の特徴が出力される。それぞれ、単語辞書「tsu」において、特徴 $f_{n2_e}$ の出現する確率が連続分布曲線222から計算され、確率値 $p_{n2_e}$ が出力される。また、単語辞書「i」において、特徴 $f_{n3_e}$ の出現する確率が連続分布曲線223から計算され、確率値 $p_{n3_e}$ が出力される。さらに、状態と状態間の状態遷移確率においても同様の処理がおこなわれ、状態遷移先を決定している。最終的に、登録された単語辞書「atsui」に対して、感情をこめて音声入力された「あつ」とい」の出現する確率値は $P_{atsui} = p_{n1_e} + p_{n2_e} + p_{n3_e}$ となる。この一連の処理を、登録された単語辞書全てにおいて計算し、計算された確率値の範囲によって感情のレベルを出力する。以上が音声の感情レベルを認識する一連の処理である。

【0041】図3は、本発明に係る他の音声および感情認識システムの各機能とその処理の流れを示すブロック図である。

【0042】音声および感情認識をおこなうために、図3に於いては、マイク301から音声が入り込まれる。取り込まれた音声であるアナログ信号は、アナログ信号をデジタル信号に変換するA/D変換器302によって、任意に決められたサンプリング周期により、アナログデータからデジタルデータに変換される。変換された音声のデジタルデータは、音声分析部303によって、雑音処理や音声分析や話者適応などの前処理がなされると共に、音声分析部303に含まれている音声パワー分析部303aで音声パワーを分析して感情のレベルが出力される。音声分析部303の出力は音声感情認識部304で処理され、音声および感情認識がなされる。ここで、音声感情認識部304で行われる音声感情認識とは、2つの処理を実行する。

【0043】第1の処理は、音声信号を解析して、それを短い時間ごとの音素として分析して、そのパターンを解析し、該当する単語や文章を辞書から選択することである。

【0044】第2の処理は、音声信号を解析して、それを短い時間(5~20ms)ごとの音素として分析して、そのパターンを解析し、話者が発声した音声の感情の度合いを示すレベルを単語や文章ごとに選択することである。



【0045】以上の2つの処理から、音声感情認識システムの出力として、音声認識結果および音声の感情レベル309を生成する。

【0046】音声感情認識部304は、音声分析部303で分析された入力音声の音声分析結果に対して、音響モデル305、単語辞書307をモデル連結部308によって連結された音素単位で照合して、単語辞書部307に登録した単語辞書307の中で、一番近い単語をピックアップする。さらに、ピックアップされた単語の入力音声を持っている感情の度合いを示すレベルを選択する。

【0047】音響モデル部305は、音声認識に用いられるモデルであり、具体的には、単語辞書部307に格納されている文字と音素との対応であり、音素の特徴が出現する確率の分布、出現した音素の特徴が次のどの特徴が現れる状態に遷移するかの確率の分布を記憶したものである。音響モデル部305は、あらかじめ声を登録しなくても、誰が話し手でもその声を認識できるいわゆる「不特定話者対応」が、一般的になってきている。このような音響モデルとしては、例えば、隠れマルコフモデル(HMM: Hidden Markov Model)を用いることができる。

【0048】単語辞書部307は、言葉、単語(名詞、動詞等)、文章を集めたものである。例えば、カーナビゲーションシステムにおいては、通り名、地名、建造物名、町名、番地、交差点名、個人住宅(個人名)、電話番号等や、必要最小限の会話に必要な言葉の集合体である。ただし、音声認識感情システムでは、特に、単語の中でも、感情を表現する単語、あるいは、感情が現われる単語で構成された単語の集合体である。より具体的には、話者が発声する「暑い」「寒い」「熱い」「冷たい」「はやく」「おそく」「大きい」「小さい」「赤い」「白い」「高く」「低く」「走れ」「進め」「戻れ」「回れ」「飛べ」等の言葉である。また、名詞等の感情表現でない単語も含まれる。この単語辞書部307に格納される単語数は、システムの能力に応じて決められるが、一つの辞書あたり、例えば、10～5000語である。

【0049】以上から、音声感情認識システム又は音声感情認識方法とは、音声信号を解析して、それを短い時間ごとの音素として分析して、そのパターンを解析し、該当する単語や文章を辞書から選択するとともに、話者が発声した音声の感情の度合いを示すレベルを単語や文章ごとに選択することである。

【0050】なお、図3に示す各処理ブロックは、複数のLSIやメモリで構成されたシステムであっても、半導体素子上に構成された一つないし複数のシステムオンチップであってもよい。また、各処理は、専用LSIや専用ICで処理するハードウェアであっても、DSPやRISCマイコン等のソフトウェアで実現したミドルウ

エアであってもよい。

【0051】図4(a)は、図3で説明した音声感情認識システムにおいて、話者が発声した音声「あつい」の音声入力波形を示すもので、横軸は時間を、縦軸は音声レベルを示す。また、図4(b)は「あつい」の音声のパワーを示したものであり、横軸に時間を、縦軸に音声のパワーを示している。

【0052】音声入力波形401は、話者が平常の音声で「あつい」と発声したときの音声波形である。音声信号は、時々刻々と変化する非定常な信号である。この音声信号を20msの短時間で切り出して見ると、定常信号と同様なスペクトル音声分析ができる。切り出された音声信号のサンプル値から、例えば、音声分析で広く用いられているLPC分析において、自己相関関数を計算すると、音声の特徴パラメータの一つとして、音声のパワーが求められる。

【0053】音声パワーを示す曲線402は、音声波形401の音声信号から計算されたパワーである。時間tに対するパワーの変化を表わしている。ここで、このパワー情報に対して、しきい値を任意に設定し、入力された音声毎にこのしきい値を超えたかどうかを観測する。この観測は、音声分析部303で行う。さらに、複数のしきい値を設け、入力された音声毎にそれぞれのしきい値を超えたかどうかを観測する。例えば、音声パワー曲線402の音声の場合は、しきい値TH1を超えているが、しきい値TH2は超えていない。すなわち、連続的にパワーが、しきい値TH1とTH2の間にある場合には、感情のレベルを1と見なし、音声分析部303は感情レベル1を出力する。

【0054】つぎに、話者が、感情を込めた強い口調の音声で「あつい」と発声したときの音声波形及び音声パワーをそれぞれ図4(c)及び図4(d)に示す。図4(c)は横軸に時間を、縦軸に音声レベルを示し、図4(d)は横軸に時間を、縦軸に音声パワーを示す。図4(c)に於いて、403は音声波形を示す。図4(d)に於いて、404は音声波形403の音声信号から計算された音声のパワーであり、時間tに対するパワーの変化を表わしている。例えば、音声パワー404の場合は、しきい値TH1を超えて、さらに、しきい値TH2を超えている。すなわち、連続的にパワーが、しきい値TH2を超えている場合には、感情のレベルを2と見なし、音声分析部303は感情レベル2を出力する。この例の場合は、感情レベルを2段階に設定したが、しきい値を増やすことによって、感情レベルをN(Nは整数)段階に設定できる。

【0055】また、音声分析部303からは、感情レベルとともに、音声認識の為の音声の特徴パラメータが時々刻々と音声感情認識部304に入力され、最終的に、音声感情認識部304からは音声認識結果「atusi」(＝あつい)と感情レベルNを示すデータ309

が出力される。

【0056】次に、図5を用いて本発明に係る音声認識システムのハードウェア構成について説明する。

【0057】音声を取り込むためのマイク501は、カーナビゲーションシステム、携帯型情報端末、PDA、ハンドヘルドPC、ゲーム、携帯型翻訳機、並びに、エアコン等の家庭電化製品等では、周囲の雑音を取り込まないために指向性をもたせた指向性マイクである。504は、マイク501により取り込まれたアナログ音声データを、デジタル音声データに変換するA/D変換器である。

【0058】音声入力用ボタン502は、音声を入力している区間を指定するためのボタンである。ボタンが押されている間、あるいは、ボタンが押された時点から音声が入力されたことをシステムに知らせる。505は、音声入力用ボタン502と、システムを接続するためのインタフェースである。

【0059】キー入力用デバイス509は、例えば、携帯型情報端末であれば、ペン入力用のデジタイザであり、ハンドヘルドPCであれば、キーボードである。また、ファミコンなどのゲーム機であれば、キャラクタ等を操作するキーパッドや、ジョイスティックである。510は、キー入力用デバイス509と、システムを接続するためのインタフェースである。

【0060】CPU503は、カーナビゲーションシステム、携帯型情報端末、PDA、ハンドヘルドPC、ゲーム、携帯型翻訳機、並びに、家庭電化製品等のメインシステムの制御と、音声感情認識システムにおける音声認識および感情認識処理を行う。図3に示す本発明の音声感情認識システムの音声分析部303、音声感情認識部304及びモデル連結部308はこのCPU503に設けられる。このCPU503には、RISCマイコンやDSPが用いられるのが、最近の潮流である。

【0061】ROM506は、音声認識用単語辞書、音響モデル、発声変形感情モデル、プログラムを格納しておく記憶装置である。また、複数の辞書や、音響モデル、発声変形感情モデルを格納しておくために、メモリカードを用いてもよい。

【0062】RAM507は、ROM506から転送された一部の辞書や、音響モデル、プログラムが格納され、また、音声感情認識処理に必要な必要最小限のワークメモリであり、ROM506に比べて、通常アクセス時間の短い半導体素子が用いられる。また、ここにはCPU503から音声認識結果及び感情のレベルを示すデータ309が入力される。

【0063】バス508は、システムにおけるデータバス、アドレスバス、制御信号バスとして用いられる。

【0064】音声感情認識結果を出力表示するためのディスプレイ512は、TFT液晶ディスプレイ等のLCDで構成し、音声認識結果および音声の感情レベルを表

示する。511は、ディスプレイ512と、システムを接続するためのインタフェースである。

【0065】音声感情認識結果を音で出力するためのスピーカ514は、音声認識結果および音声の感情レベルを音声合成して出力する。513は、音声認識結果および、音声の感情レベルをテキストから音声合成データに変換処理した後、デジタル音声合成データからアナログ音声信号に変換するD/A変換器である。

【0066】以下、本発明に係る実施形態の一例を、図6及び図8を用いて説明する。

【0067】本実施形態では、本発明の音声感情認識システムをカーエレクトロニクス製品に適用した場合について説明する。

【0068】図6(a)は本発明による音声感情認識システムをカーエレクトロニクス製品のエアコン操作に利用した場合のブロック図であり、図6(b)はこの音声感情認識システムにおける音声入力例とその認識結果を示す模式図である。

【0069】図6(a)に於いて、601は音声入力用マイク、602は音声感情認識システム、603は音声感情認識結果や、話者との双方向のやり取りを行うために、会話形式の文字情報を出力するためのディスプレイ、604は音声感情認識結果や、話者との双方向のやり取りを行うために、会話形式の文字情報を音声合成して出力するためのスピーカである。

【0070】次に、図6(b)を用いて、話者が発声した音声入力例と、音声感情認識システムが出力した認識結果例を説明する。

【0071】605は話者が音声感情システム602に対して、普通の発声で、「暑い」と発声した場合の音声感情システム602の認識結果である「暑い=感情レベル3」を示す。次に、606は話者が音声感情システムに対して、強い調子で、「暑い」と発声した時の音声感情システム602の認識結果であり、「暑い=感情レベル5」を示す。

【0072】また、607は話者が音声感情システムに対して、普通の発声で、「暑い」と発声したときの音声感情システム602の認識結果である「少し暑いですか」を示す。次に、608は話者が音声感情システムに対して、強い調子で、「暑い」と発声した時の音声感情システム602の認識結果である「かなり暑いですか」を示す。

【0073】さらに、認識結果608に対して「はい」と声感情システム602に対して発声すると、音声感情システム602その認識結果として、「車内を25℃に設定します」を出力する。実際に、車内が25℃に設定される。

【0074】また、他の本実施形態では、本発明の音声感情認識システムをファミリーコンピュータ（登録商標）等のゲーム製品に適用した場合について説明する。

【0075】図8(a)から図8(d)は、ファミリコンピュータ等のゲーム機のキャラクタの操作において、音声感情認識システムを用いた一例であり、音声によるインターフェースの音声入力例及び認識結果による動作例を示す模式図である。

【0076】図8に於いて、801、802、809及び810は、ゲーム機本体のディスプレイやゲーム機が接続されたTV等の画面である。

【0077】図8(a)に於いて、ゲームに登場するキャラクター805は、例えば、画面801に向かって左から右へと進んでいる。この操作を音声感情認識を使って行う。そこで、操作者(話者)は、発声例803に示すように「進め」と普通の音声で発声する。音声感情認識システムは、進め(susume)を認識し、さらに、感情レベルを認識する。例えば、このゲーム機における音声感情認識システムでは、感情レベルを5段階に設定したとすると、感情レベル=3と認識する。そこで、ゲーム機本体側のシステムでは、キャラクター805をキャラクター806の位置へ移動する。

【0078】図8(b)に於いて、ゲームに登場するキャラクター807は、例えば、画面802に向かって左から右へと進んでいる。この操作を音声感情認識を使って行う。そこで、操作者(話者)は、音声例804に示すように、「進め」と強い調子で発声する。音声感情認識システムは、進め(susume)を認識し、さらに、感情レベルを認識する。例えば、このゲーム機における音声感情認識システムでは、感情レベルを5段階に設定したとすると、感情レベル=5と認識する。そこで、ゲーム機本体側のシステムでは、キャラクター807はキャラクター808の位置へ大きく移動する。ここで、キャラクターの移動量は、認識された音声「進め」の感情レベルに比例する。

【0079】図8(c)に於いては、ゲームに登場するキャラクター813は、例えば、画面809に向かって左から右へと進んでいる。このとき、前方に障害物816が現われたとする。そこで、この障害物816を飛び越えなくてはならない。この操作を音声感情認識を使って行う。そこで、操作者(話者)は、音声例811に示すように「ジャンプ」と普通の音声で発声する。音声感情認識システムは、ジャンプ(jyampu)を認識し、さらに、感情レベルを認識する。例えば、このゲーム機における音声感情認識システムでは、感情レベルを5段階に設定したとすると、感情レベル=3と認識する。そこで、ゲーム機本体側のシステムでは、キャラクター813をキャラクター814の位置へ移動し、さらに、キャラクター815の位置へ移動する。

【0080】図8(d)に於いては、ゲームに登場するキャラクター817は、例えば、画面810に向かって左から右へと進んでいる。このとき、前方に障害物820が現われたとする。この障害物820は、画面80

9のときの障害物816よりも大きい。そこで、この障害物820を高く飛び越えなくてはならない。この操作を音声感情認識を使って行う。そこで、操作者(話者)は、音声例812に示すように「ジャンプ」と強い調子で発声する。音声感情認識システムは、ジャンプ(jyampu)を認識し、さらに、感情レベルを認識する。例えば、このゲーム機における音声感情認識システムでは、感情レベルを5段階に設定したとすると、感情レベル=5と認識する。そこで、ゲーム機本体側のシステムでは、キャラクター817をキャラクター818の位置へ大きく移動し、さらに、キャラクター819の位置へ移動する。キャラクターの移動量は、認識された音声「ジャンプ」の感情レベルに比例する。

【0081】

【発明の効果】本発明によれば、カーナビゲーションシステム、小型情報システム、ゲームに用いられる音声認識システムにおいて、登録した辞書の単語の文字列を音声で認識するとともに、音声認識された単語において、話者の音声を持つ感情のレベルを認識することができる音声感情認識システムを提供することができる。

【0082】また、本発明によれば、音声認識を用いたカーナビゲーションシステム、小型情報システム、ゲームにおいて、音声の感情レベルを認識できることから、限られた単語数においても、音声認識によるインタフェースのバリエーションを増やすことができ、良好な音声認識インタフェースを実現することができる。

【図面の簡単な説明】

【図1】本発明に係る音声感情認識システムの一実施例を示す示すブロック図である。

【図2】図1に示す音声感情認識システムの音響モデルおよび発声変形感情モデルを説明するための模式図である。

【図3】本発明に係る音声感情認識システムの他の実施例を示すブロック図である。

【図4】図3に示す音声感情認識システムの音声波形および音声パワーと感情レベルの関係を説明するための模式図である。

【図5】本発明のハードウェア構成を示すブロック図である。

【図6】本発明の音声認識感情システムを適用したカーナビゲーションシステムにおける音声によるインタフェースの音声入力例および認識結果例を示す模式図である。

【図7】従来の携帯型翻訳機のブロック図である。

【図8】本発明の音声認識感情システムを適用したゲーム機における音声によるインタフェースの音声入力例および認識結果による動作例を示した模式図である。

【符号の説明】

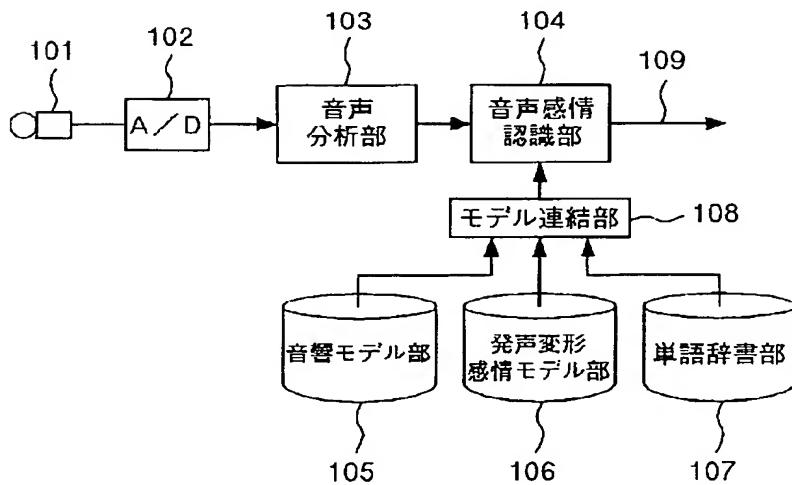
101、301、501、601…マイク、102、302、504…A/D変換器、103、303…音

声分析部、104、304…音声感情認識部、108、308…モデル連結部、105、305…音響モデル部、106…発声変形感情モデル部、107、307…単語辞書部、201…HMM音響モデル連結における「あ」の状態、202…HMM音響モデル連結における「つ」の状態、203…HMM音響モデル連結における「い」の状態、204…状態「あ」から状態「あ」へ遷移する確率、205…状態「つ」から状態「つ」へ遷移する確率、206…状態「い」から状態「い」へ遷移する確率、207…状態「あ」から状態「つ」へ遷移する確率、208…状態「つ」から状態「い」へ遷移する確率、209…状態「い」から他の状

態へ遷移する確率、211…状態「あ」の出力確率の連続分布、212…状態「つ」の出力確率の連続分布、213…状態「い」の出力確率の連続分布、214…HMM発声変形感情モデルにおける状態「あ」の出力確率の連続分布、221…HMM発声変形感情モデルにおける状態「あ」の出力確率の連続分布、222…HMM発声変形感情モデルにおける状態「つ」の出力確率の連続分布、223…HMM発声変形感情モデルにおける状態「い」の出力確率の連続分布、502…ボタン、503…CPU、506…ROM、507…RAM、509…キー、602…音声感情認識システム。

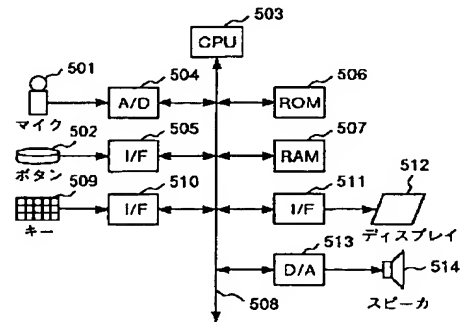
【図1】

図 1



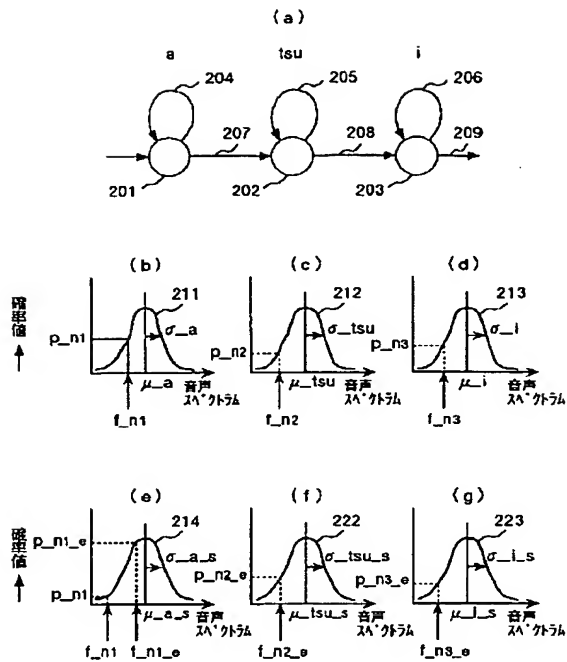
【図5】

図 5



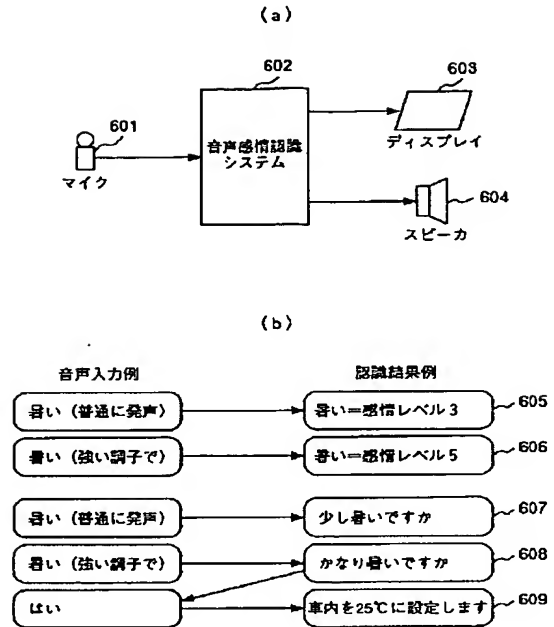
【図2】

図 2



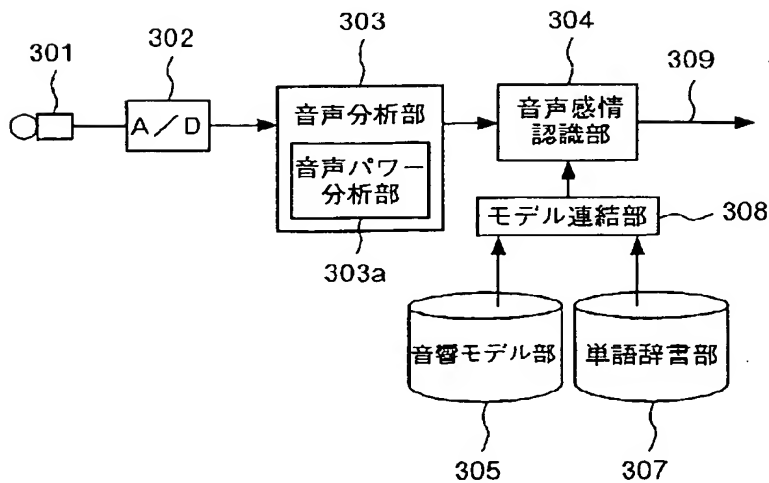
【図6】

図 6



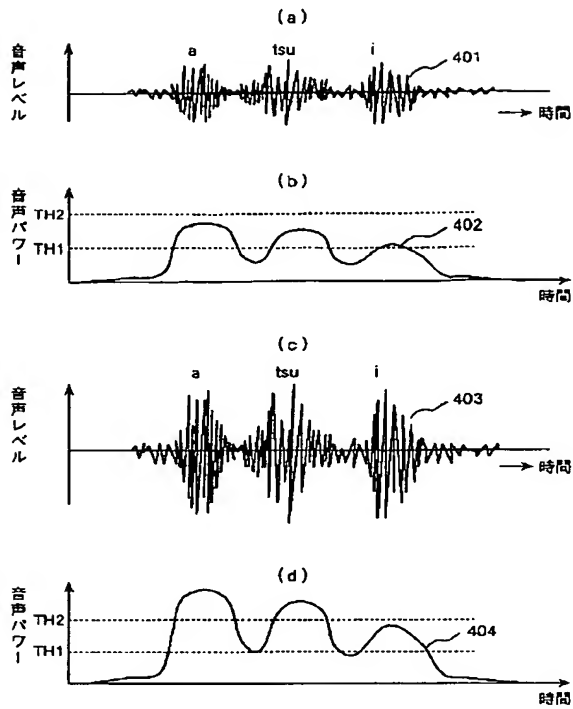
【図3】

図 3



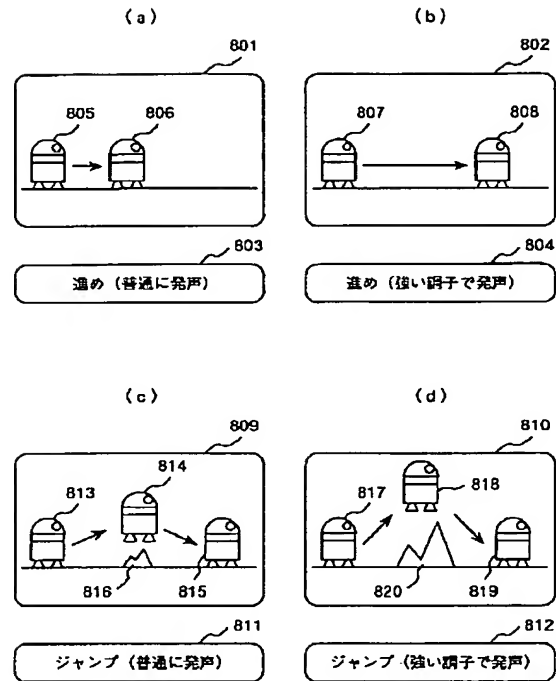
【図4】

図 4



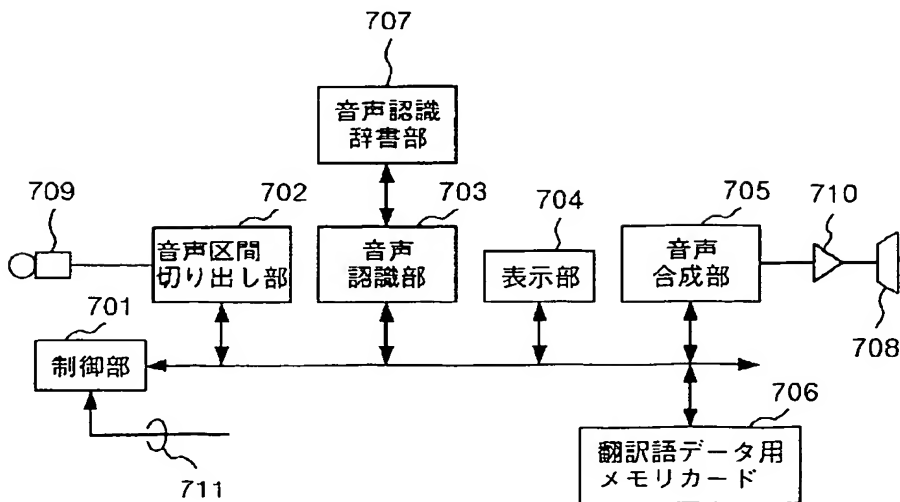
【図8】

図 8



【図7】

図 7



## フロントページの続き

(72)発明者 近藤 和夫  
東京都小平市上水本町五丁目20番1号株式  
会社日立製作所半導体事業部内  
(72)発明者 大淵 康成  
東京都国分寺市東恋ヶ窪一丁目280番地株  
式会社日立製作所中央研究所内

(72)発明者 塔下 哲司  
東京都小平市上水本町五丁目20番1号株式  
会社日立製作所半導体事業部内  
(72)発明者 石川 泰代  
東京都小平市上水本町五丁目22番1号株式  
会社日立マイコンシステム内